

Variance estimation and asymptotic confidence bands for the mean estimator of sampled functional data with high entropy unequal probability sampling designs

Hervé CARDOT^(a), Camelia GOGA^(a) and Pauline LARDIN^(a,b)

(a) Université de Bourgogne, Institut de Mathématiques de Bourgogne,

9 av. Alain Savary, 21078 DIJON, FRANCE

(b) EDF, R&D, ICAME-SOAD, 1 av. du Général de Gaulle,

92141 CLAMART, FRANCE

October 16, 2012

Abstract

For fixed size sampling designs with high entropy it is well known that the variance of the Horvitz-Thompson estimator can be approximated by the Hájek formula. The interest of this asymptotic variance approximation is that it only involves the first order inclusion probabilities of the statistical units. We extend this variance formula when the variable under study is functional and we prove, under general conditions on the regularity of the individual trajectories and the sampling design, that it asymptotically provides a uniformly consistent estimator of the variance function of the Horvitz-Thompson estimator of the mean function. Rates of convergence to the true variance function are given for rejective sampling. We deduce, under conditions on the entropy of the sampling design, that it is possible to build confidence bands whose coverage is asymptotically the desired one via simulation of Gaussian processes whose variance function is given by the Hájek formula. Finally, the accuracy of the proposed variance estimator is evaluated on samples of electricity consumption data measured every half an hour over a period of one week.

Keywords : covariance function, finite population, Hájek approximation, Horvitz-Thompson estimator, Kullback-Leibler divergence, rejective sampling, unequal probability sampling without replacement.

1 Introduction

Computing the variance of the Horvitz-Thompson estimator for unequal probability sampling designs can be difficult because the variance formula involves second order probability inclusions which are not always known. The Hájek variance formula, derived in Hájek (1964) for rejective sampling is an asymptotic approximation which only requires the knowledge of the first order inclusion probabilities and is easy to compute. It is shown in Hájek (1964) and Chen et al. (1994) that, for given first order inclusion probabilities, the rejective sampling is the fixed size sampling design with the highest entropy and the validity of this formula is closely related to the value of the entropy of the considered sampling design. Hájek (1981) proves that this approximation is also valid for the Sampford-Durbin sampling whereas Berger (1998a) gives general conditions on the relative entropy of the sampling design, also called Kullback-Leibler divergence, which entail that the use of this approximated variance formula is justified. Variants and refinements of the Hájek variance formula as well as variance estimators are proposed in Deville and Tillé (2005). Matei and Tillé (2005) show on simulations that these approximations to the variance of Horvitz-Thompson estimators work well, even for moderate sample sizes, provided that the entropy of the underlying sampling design is high enough. Recently Deville and Tillé (2005) and Fuller (2009) consider balanced, or approximately balanced, sampling algorithms which can be useful to build designs with fixed size and given inclusion probabilities by just balancing on the inclusion probabilities and relate these sampling designs with rejective sampling, so that the Hájek variance approximation remains valid.

When the aim is to build confidence intervals, the asymptotic distribution of the Horvitz-Thompson estimator is required. The Central Limit Theorem has been checked by Erdős and Rényi (1959) and Hájek (1960) for the simple random sampling without replacement, by Hájek (1964) for the rejective sampling and by Víšek (1979) for the Sampford sampling. Berger (1998b) states that the Kullback Leibler divergence of the considered sampling design, with respect to the rejective sampling, should tend to zero when the sample size gets larger for the Horvitz-Thompson estimator to be asymptotically Gaussian.

In recent studies in survey sampling the target was not a mean real value or a mean vector but a mean function (see Cardot and Josserand (2011) and Cardot et al. (2012b) for the estimation of electricity consumption curves) and one important issue was how to build confidence bands when using π ps sampling designs. A rapid technique that is well adapted for large samples has been studied in Degras (2011) and Cardot et al. (2012a). It consists in first estimating the covariance function of the mean estimator and then simulating a Gaussian process, whose covariance function is the estimated covariance function, in order to determine

the distribution of its supremum. This strategy which has been employed successfully in Cardot et al. (2012b) to build confidence bands necessitates to have an effective estimator of the variance function of the Horvitz-Thompson estimator. The aim of this work is to prove under general assumptions, on the sampling design and on the regularity of the trajectories, that the Hájek formula provides a uniformly consistent estimator of the variance function so that it is possible to assess rigorously confidence bands built with the procedure described previously.

The paper is organized as follows. The notations and our estimators are presented in Section 2. In Section 3, we state our main result, the uniform convergence of the estimated variance function under broad assumptions on the regularity of the trajectories and the sampling design. We deduce that if the Horvitz-Thompson estimator is pointwise asymptotically Gaussian, it also satisfies, under the same conditions, a functional central limit theorem and the confidence bands obtained by the Gaussian process simulation techniques have an asymptotic coverage which is the desired one. In section 4, we evaluate the performance of the covariance function estimator on samples drawn from a population of $N = 15055$ electricity consumption curves measured every half an hour over a period of one week. Note there are many ways of drawing samples with high entropy sampling distribution and given first order inclusion probabilities (see *e.g.* Brewer and Hanif (1983), Tillé (2006), Bondesson et al. (2006) and Bondesson (2010)). Because of our large population and large sample context, we use the Cube algorithm (Deville and Tillé (2004)) for which a very fast algorithm which can deal with populations of millions of units has been developed in Chauvet and Tillé (2006). The proofs are gathered in an Appendix.

2 Variance estimation and the Hájek formula

Let us consider a finite population $U_N = \{1, \dots, N\}$ of size N supposed to be known, and suppose that, for each unit k of the population U_N , we can observe a deterministic curve $Y_k = (Y_k(t))_{t \in [0, T]}$. We want to estimate the mean trajectory $\mu_N(t)$, $t \in [0, T]$, defined as follows:

$$\mu_N(t) = \frac{1}{N} \sum_{k \in U} Y_k(t). \quad (1)$$

We consider a sample s , with fixed size n , drawn from U_N according to a fixed-size sampling design $p_N(s)$, where $p_N(s)$ is the probability of drawing the sample s . The mean curve $\mu_N(t)$ is estimated by the Horvitz-Thompson estimator (Cardot et al. (2010))

$$\hat{\mu}(t) = \frac{1}{N} \sum_{k \in s} \frac{Y_k(t)}{\pi_k} = \frac{1}{N} \sum_{k \in U} \frac{Y_k(t)}{\pi_k} \mathbb{1}_k, \quad t \in [0, T], \quad (2)$$

where $\mathbb{1}_k$ is the sample membership indicator, $\mathbb{1}_k = 1$ if $k \in s$ and $\mathbb{1}_k = 0$ otherwise. We denote by $\pi_k = \mathbb{E}_p(\mathbb{1}_k)$ the first order inclusion probability of unit k with respect to the sampling design $p_N(s)$ and we suppose that $\pi_k > 0$, for all units k in U . It is well known that, for each value of $t \in [0, T]$, $\hat{\mu}(t)$ is a design-unbiased estimator of $\mu_N(t)$, *i.e.* $\mathbb{E}_p(\hat{\mu}(t)) = \mu_N(t)$. We denote by $\pi_{kl} = \mathbb{E}_p(\mathbb{1}_k \mathbb{1}_l)$ the second order inclusion probabilities and we suppose that $\pi_{kl} > 0$.

Since the sample size is fixed, the variance $\gamma_p(t, t)$ for each instant t of the estimator $\hat{\mu}(t)$ is given by the Yates and Grundy formula (see Yates and Grundy (1953) and Sen (1953)),

$$\gamma_p(t, t) = -\frac{1}{2} \frac{1}{N^2} \sum_{k \in U} \sum_{l \in U, l \neq k} (\pi_{kl} - \pi_k \pi_l) \left(\frac{Y_k(t)}{\pi_k} - \frac{Y_l(t)}{\pi_l} \right)^2, \quad (3)$$

and it is straightforward to express the covariance $\gamma_p(r, t)$ of $\hat{\mu}$ between two instants r and t , as follows

$$\gamma_p(r, t) = -\frac{1}{2} \frac{1}{N^2} \sum_{k \in U} \sum_{l \in U, l \neq k} (\pi_{kl} - \pi_k \pi_l) \left(\frac{Y_k(r)}{\pi_k} - \frac{Y_l(r)}{\pi_l} \right) \left(\frac{Y_k(t)}{\pi_k} - \frac{Y_l(t)}{\pi_l} \right). \quad (4)$$

The variance formula (3) clearly indicates that if the first order inclusion probabilities are chosen to be approximately proportional to $Y_k(t)$, the variance of the estimator $\hat{\mu}(t)$ will be small. Thus, if we have an auxiliary variable X , whose value x_k , supposed to be positive, is known for all the units $k \in U$ and if X is correlated with the variable of interest, it can be interesting to consider a sampling design whose first order inclusion probabilities are given by

$$\pi_k = n \frac{x_k}{\sum_U x_k}.$$

There are many ways of building sampling designs with given first order inclusion probabilities (see *e.g.* Brewer and Hanif (1983) and Tillé (2006)) and we focus here on the designs with high entropy, where the entropy of a sampling design p_N (a discrete probability distribution on U_N) is defined by

$$H(p_N) = - \sum_s p_N(s) \ln(p_N(s))$$

with the convention $0 \ln 0 = 0$. Chen et al. (1994) have proved that for given first order inclusion probabilities, the rejective sampling, or conditional Poisson sampling, is the fixed size sampling design with the highest entropy. Then, a key result is the following uniform approximation to the second order inclusion probabilities,

$$\pi_{kl} = \pi_k \pi_l \left\{ 1 - \frac{(1 - \pi_k)(1 - \pi_l)}{d(\pi)} [1 + o(1)] \right\} \quad (5)$$

where $d(\pi) = \sum_U \pi_k(1 - \pi_k)$ is supposed to tend to infinity. This approximation, which is satisfied for the rejective sampling and the Sampford-Durbin sampling (see Hájek (1981)),

appears to be very efficient when the sample size is large enough (and thus the value of d is also large) and the entropy of the sampling design is close to the maximum entropy. It also ensures that the variance estimator given below is always positive.

By plugging in approximation (5) in (4), we obtain the Hájek approximation $\gamma_H(r, t)$ of the covariance $cov(\hat{\mu}(t), \hat{\mu}(r))$ in the functional case : for all $(r, t) \in [0, T]^2$,

$$\gamma_H(r, t) = \frac{1}{N^2} \left[\sum_U \frac{Y_k(t)Y_k(r)}{\pi_k} (1 - \pi_k) - \frac{1}{d(\pi)} \sum_U \sum_U (1 - \pi_k)(1 - \pi_l) Y_k(t)Y_l(r) \right]. \quad (6)$$

We consider in the following two variance estimators

$$\hat{\gamma}_H(r, t) = \frac{1}{N^2} \frac{\hat{d}(\pi)}{d(\pi)} \left[\sum_s \frac{1 - \pi_k}{\pi_k^2} Y_k(t)Y_k(r) - \frac{1}{\hat{d}(\pi)} \sum_s \sum_s \frac{1 - \pi_k}{\pi_k} \frac{1 - \pi_l}{\pi_l} Y_k(t)Y_l(r) \right], \quad (7)$$

and

$$\hat{\gamma}_H^*(r, t) = \frac{1}{N^2} \left[\sum_s \frac{1 - \pi_k}{\pi_k^2} Y_k(t)Y_k(r) - \frac{1}{\hat{d}(\pi)} \sum_s \sum_s \frac{1 - \pi_k}{\pi_k} \frac{1 - \pi_l}{\pi_l} Y_k(t)Y_l(r) \right] \quad (8)$$

where $\hat{d}(\pi) = \sum_s (1 - \pi_k)$. Note that $\hat{\gamma}_H(r, t)$ is the functional analogue of the slightly modified variance estimator proposed by Berger (1998a) in the real case. More exactly, the variance estimator considered by Berger (1998a) is $\hat{\gamma}_H(r, t)$ multiplied by the correction factor $n/(n-1)$ that allows to obtain the exact expression for simple random sampling without replacement. The second estimator, $\hat{\gamma}_H^*(r, t)$ is the extension to the functional case of the Deville and Tillé (2005)'s estimator and it has been successfully used in a very recent study by Cardot et al. (2012b) to build confidence bands for the mean electricity consumption curve.

We can easily show the following property.

Proposition 2.1. *If, for all $t \in [0, T]$, there is a constant c_t such that $Y_k(t) = c_t \pi_k$ then $\gamma_H(r, t) = 0$ and $\hat{\gamma}_H(r, t) = \hat{\gamma}_H^*(r, t) = 0$.*

With real data, we do not observe $Y_k(t)$ at all instants t in $[0, T]$ but only for a finite set of D measurement times, $0 = t_1 < \dots < t_D = T$. In functional data analysis, when the noise level is low and the grid of discretization points is fine, it is usual to perform a linear interpolation or a smoothing of the discretized trajectories in order to obtain approximations of the trajectories at every instant t (cf. Ramsay and Silverman (2005)). When there are no measurement errors and when the trajectories are regular enough, Cardot and Josserand (2011) showed that linear interpolation can provide sufficiently accurate approximations of the trajectories. Thus, for each unit k in the sample s , we build the interpolated trajectory

$$Y_{k,d}(t) = Y_k(t_i) + \frac{Y_k(t_{i+1}) - Y_k(t_i)}{t_{i+1} - t_i} (t - t_i), \quad t \in [t_i, t_{i+1}],$$

and define the estimator of the mean curve $\mu_N(t)$ based on the discretized observations as follows

$$\hat{\mu}_d(t) = \frac{1}{N} \sum_s \frac{Y_{k,d}(t)}{\pi_k}, \quad t \in [t_i, t_{i+1}]. \quad (9)$$

Its covariance is then estimated by

$$\hat{\gamma}_{H,d}(r, t) = \frac{1}{N^2} \frac{\hat{d}(\pi)}{d(\pi)} \left[\sum_s \frac{1 - \pi_k}{\pi_k^2} Y_{k,d}(t) Y_{k,d}(r) - \frac{1}{\hat{d}(\pi)} \sum_s \sum_s \frac{1 - \pi_k}{\pi_k} \frac{1 - \pi_l}{\pi_l} Y_{k,d}(t) Y_{l,d}(r) \right], \quad (10)$$

and we show in the next section that it is an uniformly consistent estimator of the variance function. Replacing $Y_k(t)$ by $Y_{k,d}(t)$ in (8), yields the variance estimator $\hat{\gamma}_{H,d}^*(r, t)$ based on discretized values.

3 Asymptotic properties

All the proof are postponed in an Appendix.

3.1 Assumptions

To demonstrate the asymptotic properties, we must suppose that the sample size and population size become large. Therefore, we adopt the asymptotic approach of Hájek (1964), assuming that $d(\pi) \rightarrow \infty$. Note that this assumption implies that $n \rightarrow \infty$ and $N - n \rightarrow \infty$. We consider a sequence of growing and nested populations U_N with size N tending to infinity and a sequence of samples s_N of size n_N drawn from U_N according to the sampling design $p_N(s_N)$. The first and second order inclusion probabilities are respectively denoted by π_{kN} and π_{klN} . For simplicity of notations and when there is no ambiguity, we drop the subscript N . To prove our asymptotic results we need to introduce the following assumptions.

A1. We assume that $\lim_{N \rightarrow \infty} \frac{n}{N} = \pi \in]0, 1[$.

A2. We assume that $\min_{k \in U} \pi_k \geq \lambda > 0$, $\min_{k \neq l} \pi_{kl} \geq \lambda^* > 0$ and

$$\pi_{kl} = \pi_k \pi_l \left\{ 1 - \frac{(1 - \pi_k)(1 - \pi_l)}{d(\pi)} [1 + o(1)] \right\}.$$

A3. There are two positive constants C_2 and C_3 and $\beta > 1/2$ such that, for all N and for all $(r, t) \in [0, T] \times [0, T]$,

$$\frac{1}{N} \sum_{k \in U} (Y_k(0))^2 < C_2 \quad \text{and} \quad \frac{1}{N} \sum_{k \in U} (Y_k(t) - Y_k(r))^2 < C_3 |t - r|^{2\beta}.$$

A4. There are two positive constants C_4 and C_5 such that, for all N and for all $(r, t) \in [0, T] \times [0, T]$,

$$\frac{1}{N} \sum_{k \in U} (Y_k(0))^4 < C_4 \quad \text{and} \quad \frac{1}{N} \sum_{k \in U} (Y_k(t) - Y_k(r))^4 < C_5 |t - r|^{4\beta}.$$

A5. We assume that

$$\lim_{N \rightarrow \infty} \max_{(k_1, l_1, k_2, l_2) \in D_{4,N}} |\mathbb{E}_p [(\mathbb{1}_{k_1 l_1} - \pi_{k_1} \pi_{l_1})(\mathbb{1}_{k_2 l_2} - \pi_{k_2} \pi_{l_2})]| \rightarrow 0$$

where $\mathbb{1}_{kl}$ is the sample membership of the couple (k, l) and $D_{4,N}$ is the set of all distinct quadruples (i_1, \dots, i_4) from U_N .

Assumptions **A1** and **A2** are classical hypotheses in survey sampling and deal with the first and second order inclusion probabilities. They are satisfied for many usual sampling designs with fixed size (see for example Hájek (1981)). They directly imply that $cn \leq d(\pi) \leq n$, for some strictly positive constant c . The assumption **A2** implies that $\limsup_{N \rightarrow \infty} n \max_{k \neq l} |\pi_{kl} - \pi_k \pi_l| < C_1 < \infty$. It also ensures that the Yates-Grundy variance estimator is always positive since $\pi_{kl} \leq \pi_k \pi_l$.

Assumption **A3** and **A4** are regularity conditions on the individual trajectories. Even if point-wise consistency, for each fixed value of t , can be proved without any condition on β , these regularity conditions are required to get uniform convergence of the mean estimator (see Cardot and Josserand (2011)). Note finally that assumption **A5** is true for SRSWOR, stratified sampling and rejective sampling (Boistard et al. (2012)). More generally, it also holds for unequal probability designs with large entropy as shown in the following proposition. Let us recall before the definition of the Kullback-Leibler divergence $K(p_N, p_{rej})$,

$$K(p_N, p_{rej}) = \sum_s p_N(s) \ln \left(\frac{p_N(s)}{p_{rej}(s)} \right), \quad (11)$$

which measures how a sampling distribution $p_N(s)$ is distant from a reference sampling distribution, chosen here to be the rejective sampling $p_{rej}(s)$ since it is the design with maximum entropy for given first order inclusion probabilities. We can now state the following proposition which gives an upper bound for the rates of convergence to zero of the quantity in **A4** in terms of Kullback-Leibler divergence with respect to the rejective sampling. We consider a sampling design p_N with the same first order inclusion probabilities as p_{rej} . We have

Proposition 3.1. *If $d(\pi) \rightarrow \infty$, then*

$$\max_{(k_1, l_1, k_2, l_2) \in D_{4,N}} |\mathbb{E}_p [(\mathbb{1}_{k_1 l_1} - \pi_{k_1} \pi_{l_1})(\mathbb{1}_{k_2 l_2} - \pi_{k_2} \pi_{l_2})]| \leq \frac{C}{d(\pi)} + \sqrt{K(p_N, p_{rej})/2}$$

for some constant C .

A direct consequence of Proposition 3.1 is that assumption **A5** is satisfied for rejective sampling as well as for the Sampford-Durbin design, whose Kullback-Leibler divergence, with respect to the rejective sampling, tends to zero as the sample size n tends to infinity (see Berger (1998b)). Note also that the Kullback-Leibler divergence has been approximated asymptotically for other sampling designs such as Pareto sampling in Lundqvist (2007).

3.2 Convergence of the estimated variance

Let us first recall Proposition 3.3 in Cardot and Josserand (2011) which states that the estimator $\hat{\mu}_d$ is asymptotically design unbiased and uniform consistent under mild assumption. More precisely, if assumptions (A1)-(A3) hold and if the discretization scheme satisfies $\lim_{N \rightarrow \infty} \max_{i=\{1, \dots, d_N-1\}} |t_{i+1} - t_i|^{2\beta} = o(n^{-1})$, then for some constant C

$$\sqrt{n} \mathbb{E}_p \left\{ \sup_{t \in [0, T]} |\hat{\mu}_d(t) - \mu_N(t)| \right\} \leq C.$$

We can now state our first result which indicates that $\hat{\gamma}_{H,d}$ is consistent pointwise and the variance function estimator of the estimated mean trajectory is uniformly consistent. Note that additional assumptions on the sampling design are required to obtain the rate of convergence.

Proposition 3.2. *Assume (A1)-(A5) hold and the sequence of discretization schemes satisfies $\lim_{N \rightarrow \infty} \max_{i=\{1, \dots, d_N-1\}} |t_{i+1} - t_i| = o(1)$. When N tends to infinity,*

$$n \mathbb{E}_p \{ | \hat{\gamma}_{H,d}(r, t) - \gamma_p(r, t) | \} \rightarrow 0$$

for all $(r, t) \in [0, T]^2$ and

$$n \mathbb{E}_p \left\{ \sup_{t \in [0, T]} | \hat{\gamma}_{H,d}(t, t) - \gamma_p(t, t) | \right\} \rightarrow 0.$$

We can state the same result for the second variance estimator, $\hat{\gamma}_{H,d}^*$.

Proposition 3.3. *Assume (A1)-(A5) hold and the sequence of discretization schemes satisfies $\lim_{N \rightarrow \infty} \max_{i=\{1, \dots, d_N-1\}} |t_{i+1} - t_i| = o(1)$. When N tends to infinity,*

$$n \mathbb{E}_p \{ | \hat{\gamma}_{H,d}^*(r, t) - \gamma_p(r, t) | \} \rightarrow 0$$

for all $(r, t) \in [0, T]^2$ and

$$n \mathbb{E}_p \left\{ \sup_{t \in [0, T]} | \hat{\gamma}_{H,d}^*(t, t) - \gamma_p(t, t) | \right\} \rightarrow 0.$$

A sharper result can be stated for the particular case of rejective sampling for which accurate approximations to the multiple inclusion probabilities are available (see Boistard et al. (2012)).

Proposition 3.4. *Suppose that the sample is selected with the rejective sampling design. Assume (A1)-(A4) hold and the sequence of discretization schemes satisfies $\lim_{N \rightarrow \infty} \max_{i=\{1, \dots, d_N-1\}} |t_{i+1} - t_i|^{2\beta} = O(n^{-1})$. Then, for all $(r, t) \in [0, T]^2$*

$$n^3 \mathbb{E}_p \left[(\hat{\gamma}_{H,d}(r, t) - \gamma_p(r, t))^2 \right] \leq C$$

for some positive constant C .

We can note in the proof, given in the Appendix, that the approximation error to the true variance by the Hájek formula is asymptotically negligible compared to the sampling error.

3.3 Asymptotic normality and confidence bands

Let us assume that the Horvitz-thompson estimator satisfies a Central Limit Theorem for real valued quantities with new moment conditions

A6. There is some $\delta > 0$, such that $N^{-1} \sum_{k \in U_N} |Y_k(t)|^{2+\delta} < \infty$ for all $t \in [0, T]$, and $\{\gamma_p(t, t)\}^{-1/2} \{\hat{\mu}(t) - \mu(t)\} \rightarrow \mathcal{N}(0, 1)$ in distribution when N tends to infinity.

Cardot and Josserand (2011) have shown that under the previous assumptions, the central limit theorem also holds in the space of continuous functions. More precisely, if assumptions (A1)-(A3) and (A6) hold and the discretization points satisfy

$\lim_{N \rightarrow \infty} \max_{i=\{1, \dots, d_N-1\}} |t_{i+1} - t_i|^{2\beta} = o(n^{-1})$, we have

$$\sqrt{n}(\hat{\mu}_d - \mu) \rightarrow Z \text{ in distribution in } C[0, T]$$

where Z is a Gaussian random function taking values in $C[0, T]$ with mean 0 and covariance function $\gamma_Z(r, t) = \lim_{N \rightarrow \infty} n\gamma_{p_N}(r, t)$. This important result gives a theoretical justification of the confidence bands built as follows. We examine now the asymptotic coverage of confidence bands for μ_N of the form

$$\left\{ \left[\hat{\mu}_d(t) \pm c \frac{\hat{\sigma}(t)}{\sqrt{n}} \right], t \in [0, T] \right\}, \quad (12)$$

where c is a suitable number and $\hat{\sigma}(t) = \sqrt{n\hat{\gamma}_{H,d}(t, t)}$.

Given a confidence level $1 - \alpha \in]0, 1[$, one way to build such confidence bands, that is to say one way to find an adequate value for c_α , is to perform simulations of centered Gaussian functions \hat{Z} defined on $[0, T]$ with mean 0 and covariance function $n\hat{\gamma}_{H,d}(r, t)$ and then compute the quantile of order $1 - \alpha$ of $\sup_{t \in [0, T]} |\hat{Z}(t)/\hat{\sigma}(t)|$. In other words, we look for

a constant c_α , which is in fact a random variable since it depends on the estimated covariance function $\hat{\gamma}_{MA,d}$, such that

$$\mathbb{P} \left(|\hat{Z}(t)| \leq c_\alpha \frac{\hat{\sigma}(t)}{\sqrt{n}}, \forall t \in [0, T] \mid \hat{\gamma}_{H,d} \right) = 1 - \alpha. \quad (13)$$

Next proposition provides a rigorous justification of this latter technique :

Proposition 3.5. *Assume (A1)-(A6) hold and the discretization scheme satisfies*

$$\lim_{N \rightarrow \infty} \max_{i=\{1, \dots, d_N-1\}} |t_{i+1} - t_i|^{2\beta} = o(n^{-1}).$$

Let Z be a Gaussian process with mean zero and covariance function γ_Z . Let (\hat{Z}_N) be a sequence of processes such that for each N , conditionally on $\hat{\gamma}_{H,d}$ defined in (10), \hat{Z}_N is Gaussian with mean zero and covariance $n\hat{\gamma}_{H,d}$. Then for all $c > 0$, as $N \rightarrow \infty$, the following convergence holds in probability:

$$\mathbb{P} \left(|\hat{Z}_N(t)| \leq c \hat{\sigma}(t), \forall t \in [0, T] \mid \hat{\gamma}_{H,d} \right) \rightarrow \mathbb{P} (|Z(t)| \leq c \sigma(t), \forall t \in [0, T]),$$

where $\hat{\sigma}(t) = \sqrt{n\hat{\gamma}_{H,d}(t, t)}$ and $\sigma(t) = \sqrt{\gamma_Z(t, t)}$.

The proof of Proposition 3.5 is very similar to the proof of Proposition 3.5 in Cardot et al. (2012c) and is thus omitted. As in Cardot et al. (2012a), it is possible to deduce from previous proposition that the chosen value $\hat{c}_\alpha = c_\alpha(\hat{\gamma}_{H,d})$ provides asymptotically the desired coverage since it satisfies

$$\lim_{N \rightarrow \infty} \mathbb{P} \left(\mu(t) \in \left[\hat{\mu}_d(t) \pm \hat{c}_\alpha \frac{\hat{\sigma}(t)}{\sqrt{n}} \right], \forall t \in [0, T] \right) = 1 - \alpha.$$

4 Example: variance estimation for electricity consumption curves

In this section, we evaluate the performance of the estimators $\hat{\gamma}_{H,d}^*(r, t)$ and $\hat{\gamma}_{H,d}(r, t)$ of the functional variance $\gamma_p(r, t)$ of $\hat{\mu}_d(t)$. Simulation studies not reported here showed that the estimators $\hat{\gamma}_{H,d}^*(r, t)$ and $\hat{\gamma}_{H,d}(r, t)$ conduct very similarly asymptotically. This is why we only give below the simulation results for $\hat{\gamma}_{H,d}(r, t)$.

We use the same data frame as in Cardot et al. (2012b). More exactly, we have a population U of $N = 15055$ electricity consumption curves measured every half an hour during one week, so that there are $\mathcal{D} = 336$ time points. The mean consumption during the previous week for each meter k , denoted x_k , is used as an auxiliary variable. This variable is strongly correlated to the consumption curve $Y_k(t)$ (the pointwise correlation is always larger than 0.80) and is inexpensive to transmit.

We select samples s of size n drawn with inclusion probabilities π_k proportional to the past mean electricity consumption. This means that $\pi_k = n \frac{x_k}{\sum_U x_k}$. As mentioned in Deville and

Tillé (2005), this kind of sampling may be viewed as a balanced sampling with the balancing variable $\boldsymbol{\pi} = (\pi_1, \dots, \pi_N)$. The sample was drawn using the fast version (see Chauvet and Tillé (2006)) of the cube algorithm (see Deville and Tillé (2004)). As suggested in Chauvet (2007), a random sort of the population is made before the sample selection. The true mean consumption curve observed in the population U and one estimation obtained from a sample s' of size $n = 1500$ are drawn in Figure 1.

The inclusion probabilities π_{kl} being unknown, an empirical estimation of the covariance function γ_p is given, from $J = 10000$ simulations, by

$$\gamma_{emp}(r, t) = \frac{1}{J-1} \sum_{j=1}^J (\hat{\mu}_{d,j}(t) - \hat{\bar{\mu}}_d(t))(\hat{\mu}_{d,j}(r) - \hat{\bar{\mu}}_d(r)) \quad (14)$$

with $\hat{\mu}_{d,j}(t) = \frac{1}{N} \sum_{k \in s_j} \frac{Y_{k,d}(t)}{\pi_k}$, $\hat{\bar{\mu}}_d(t) = \frac{1}{J} \sum_{j=1}^J \hat{\mu}_{d,j}(t)$ and $(r, t) \in [0, T]$. The empirical variance function γ_{emp} (solid line) of estimator $\hat{\mu}_d$, the Hájek approximation γ_H (dotted line) and one estimation $\hat{\gamma}_{H,d}$ (dashed line) obtained from the same sample s' are drawn in Figure 2.

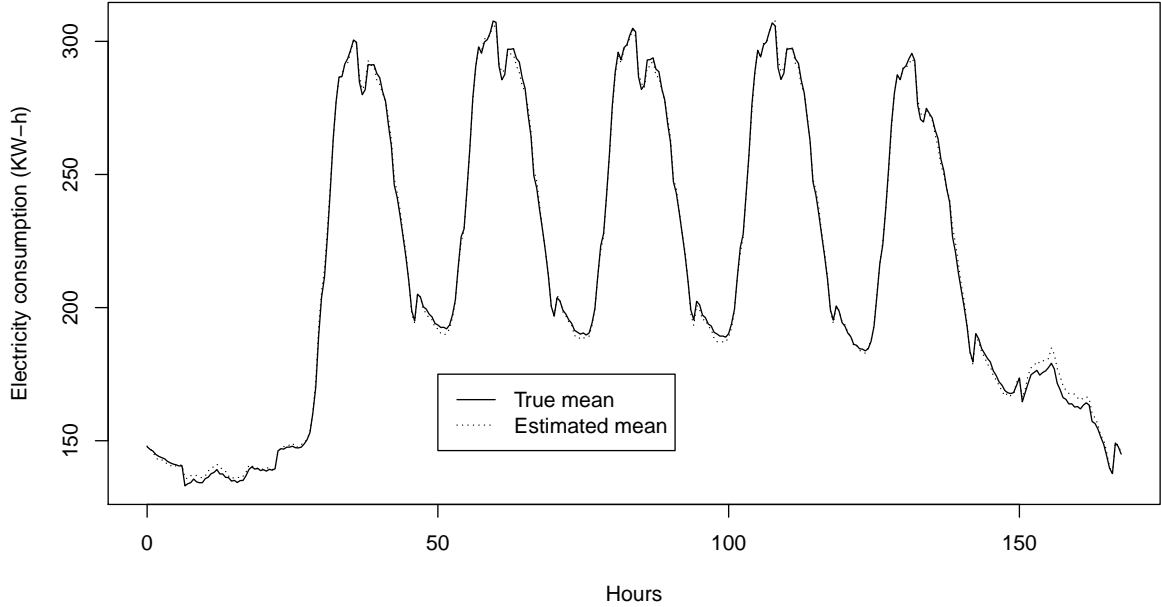


Figure 1: Mean consumption curve and its Horvitz-Thompson estimation obtained from sample s' , with $n = 1500$.

To evaluate the performance of estimator $\hat{\gamma}_{H,d}$, we consider different sample sizes, $n = 250$, $n = 500$ and $n = 1500$. The corresponding values of $d(\pi)$ are $d(\pi) = 241.2$, $d(\pi) = 464.7$ and $d(\pi) = 1202.3$.

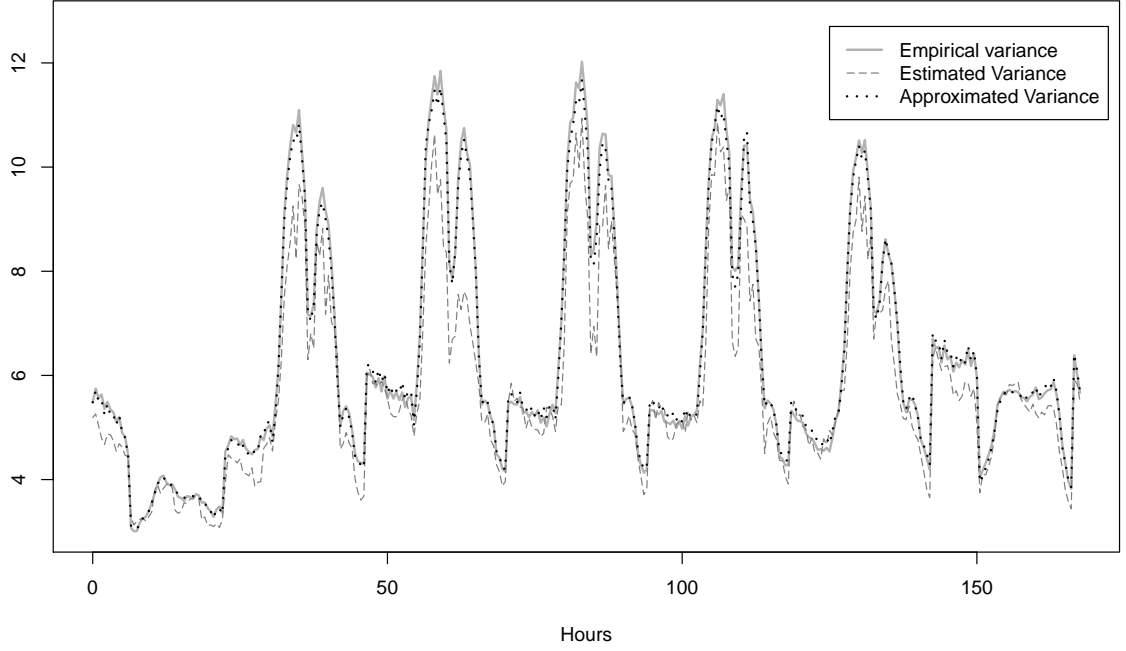


Figure 2: Empirical variance γ_{emp} (solid line), Hájek's approximation γ_H (dotted line) and variance estimation $\hat{\gamma}_{H,d}$ (dashed line) obtained from sample s' , with $n = 1500$.

For each sample size, we draw $I = 10000$ samples and we compute the following quadratic loss criterion

$$R(\hat{\gamma}_{H,d}) = \frac{1}{\mathcal{D}} \sum_{\mathbf{d}=1}^{\mathcal{D}} \frac{|\hat{\gamma}_H(t_{\mathbf{d}}, t_{\mathbf{d}}) - \gamma_{emp}(t_{\mathbf{d}}, t_{\mathbf{d}})|^2}{\gamma_{emp}(t_{\mathbf{d}}, t_{\mathbf{d}})^2} \simeq \int \frac{|\hat{\gamma}_H(t, t) - \gamma_{emp}(t, t)|^2}{\gamma_{emp}(t, t)^2} dt. \quad (15)$$

We also compute the relative mean squared error,

$$RMSE = \frac{1}{I} \sum_{i=1}^I R^{(i)}(\hat{\gamma}_{H,d}) = RB(\hat{\gamma}_{H,d})^2 + RV(\hat{\gamma}_{H,d}), \quad (16)$$

where $R^{(i)}(\hat{\gamma}_{H,d})$ is the value of $R(\hat{\gamma}_{H,d})$ computed for the i th simulation; $RB(\hat{\gamma}_{H,d})$ and $RV(\hat{\gamma}_{H,d})$ are the relative bias, respectively the relative variance of estimator $\hat{\gamma}_{H,d}$. Note that $RB(\hat{\gamma}_{H,d})^2$ is given by

$$RB(\hat{\gamma}_{H,d})^2 = \frac{1}{\mathcal{D}} \sum_{\mathbf{d}=1}^{\mathcal{D}} \left(\frac{\bar{\hat{\gamma}}_{H,d}(t_{\mathbf{d}}, t_{\mathbf{d}}) - \gamma_{emp}(t_{\mathbf{d}}, t_{\mathbf{d}})}{\gamma_{emp}(t_{\mathbf{d}}, t_{\mathbf{d}})} \right)^2$$

where $\bar{\hat{\gamma}}_{H,d}(t_{\mathbf{d}}, t_{\mathbf{d}}) = \sum_{i=1}^I \hat{\gamma}_{H,d}^{(i)}(t_{\mathbf{d}}, t_{\mathbf{d}}) / I$ and $\hat{\gamma}_{H,d}^{(i)}(t_{\mathbf{d}}, t_{\mathbf{d}})$ is the variance estimation obtained with the i th simulation.

Sample Size	$RMSE$	$RB(\hat{\gamma}_{H,d})^2$	$R(\hat{\gamma}_{H,d})$				
			5%	1 st quartile	median	3 rd quartile	95%
250	0.9473	0.0004	0.0188	0.0298	0.0446	0.0748	0.4326
500	0.3428	0.0002	0.0121	0.0191	0.0278	0.0456	0.3510
1500	0.1406	0.0003	0.006	0.0097	0.0144	0.0272	0.0929

Table 1: $RMSE$, $RB(\hat{\gamma}_{H,d})^2$ and estimation errors according to criterion $R(\hat{\gamma}_{H,d})$ for different sample sizes, with $I = 10000$ simulations.

The estimation errors are presented in Table 1 for the three considered sample sizes. We first note that the values of relative bias $RB(\hat{\gamma}_{H,d})$ are very low, meaning that the Hájek's formula provides, in our relatively large sample context, a very good approximation to the variance. The median error for $R(\hat{\gamma}_{H,d})$ is slightly larger but remains small (always less than 5%), even for moderate sample sizes ($n=250$). This means that the most important part of the variance estimation error is due to the sampling error. We have drawn in Figure 3 the approximation error $\gamma_{emp}(t, r) - \gamma_{H,d}(t, r)$ and in Figure 4 the estimation error $\gamma_{emp}(t, r) - \hat{\gamma}_{H,d}(t, r)$ for $t, r \in \{1, \dots, \mathcal{D}\}$, corresponding to a sample of size $n = 1500$ with an estimation error close to the median value of the global risk, $R(\hat{\gamma}_{H,d}) = 0.0144$. It appears that the largest estimation errors for the variance occur when the level of consumption is high. We can also observe in these Figures a kind of periodic pattern which can be related to the daily electricity consumption behavior.

Nevertheless, we also note that the relative mean squared error $RMSE$, which is approximately equal to the relative variance of the estimator $\hat{\gamma}_{H,d}$, is rather high, especially for small sample sizes ($n = 250$). Looking at the 95 % quantiles of $R(\hat{\gamma}_{H,d})$ in Table 1, we can deduce that bad variance estimations only occur in rare cases but with very large errors. We can note in Figure 5, which represents the distribution of the sampling weights at a logarithmic scale, that there are many large outlying values, especially when the sample size is not very large. The bad performance of the variance estimator, in terms of RMSE, is in fact due to a few individuals in the population that have both a very small inclusion probability π_k and a consumption level Y_k that can be very high at some instants of the period. Their selection in the sample, which occurs rarely, leads to an overestimation of the mean curve and to a large error $R(\hat{\gamma}_{H,d})$ when estimating the variance at these instants. One possible way to deal with this issue and that will be explored in a another work would consist in correcting the sampling weights of the most influential units of the sample (see *e.g* Beaumont and Rivest (2009)) in order to get a more stable variance estimator.

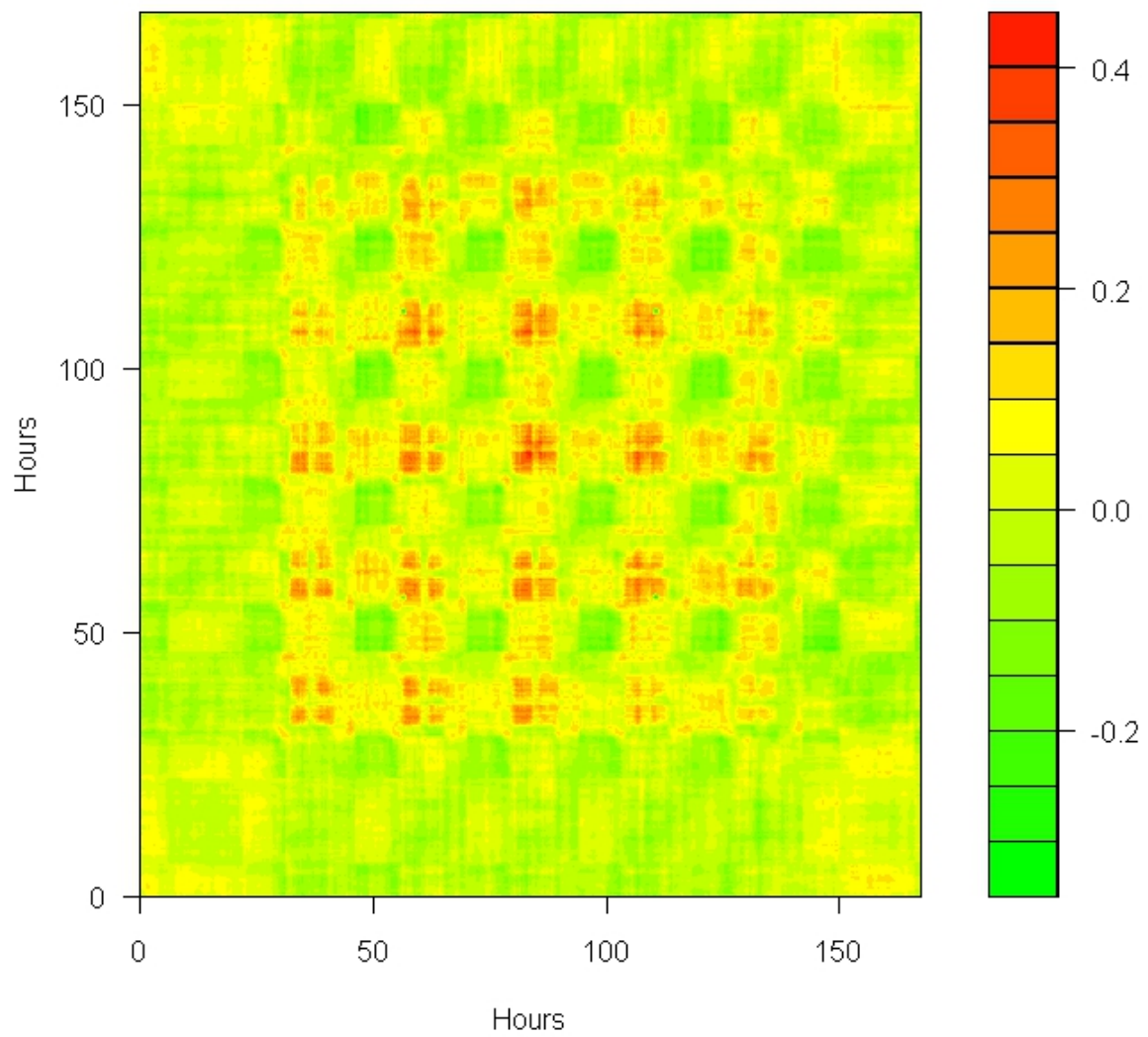


Figure 3: Approximation error $\gamma_{emp} - \gamma_{H,d}$ for a sample of size $n = 1500$.

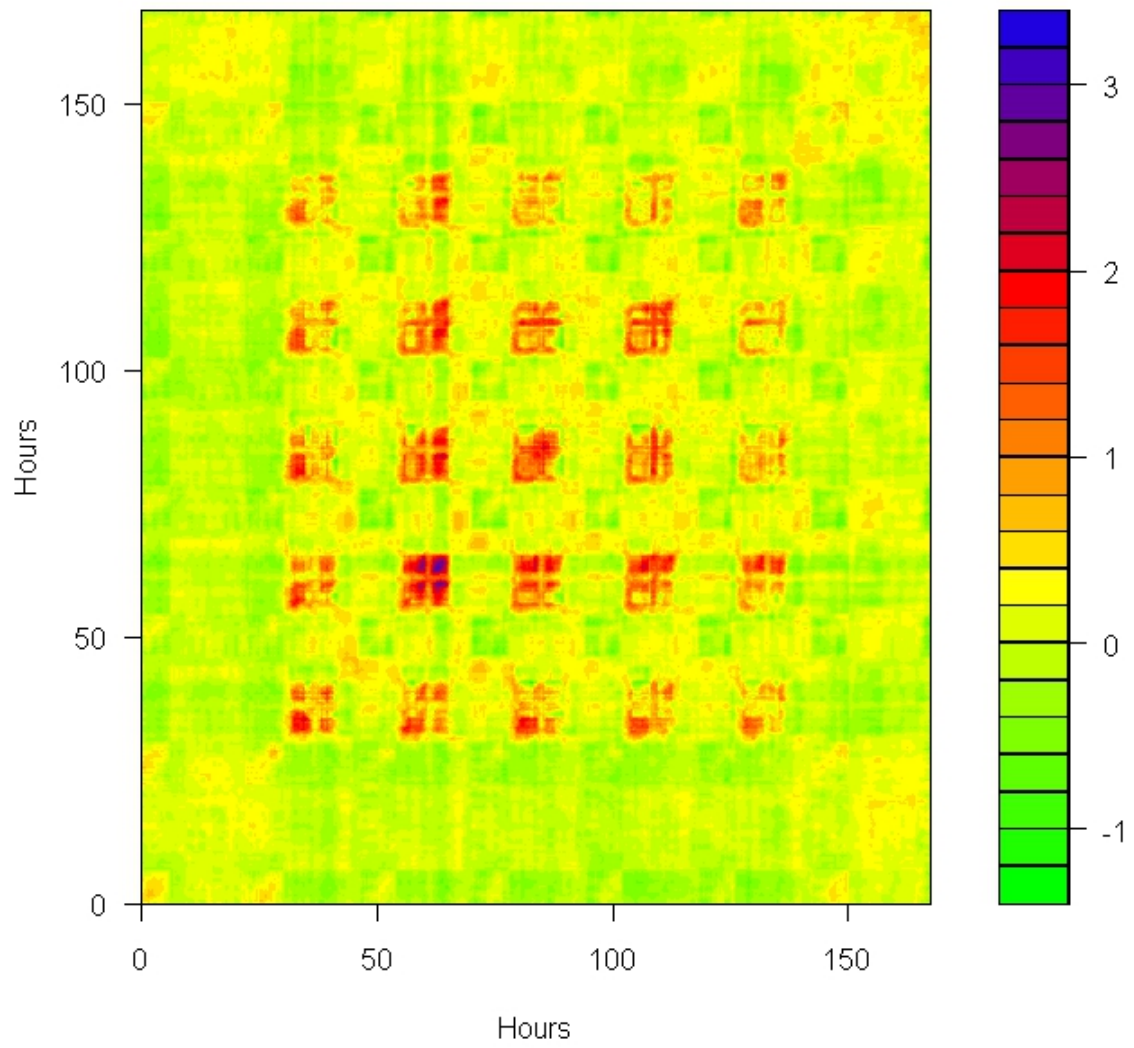


Figure 4: Estimation error $\gamma_{emp} - \hat{\gamma}_{H,d}$ for a sample of size $n = 1500$.

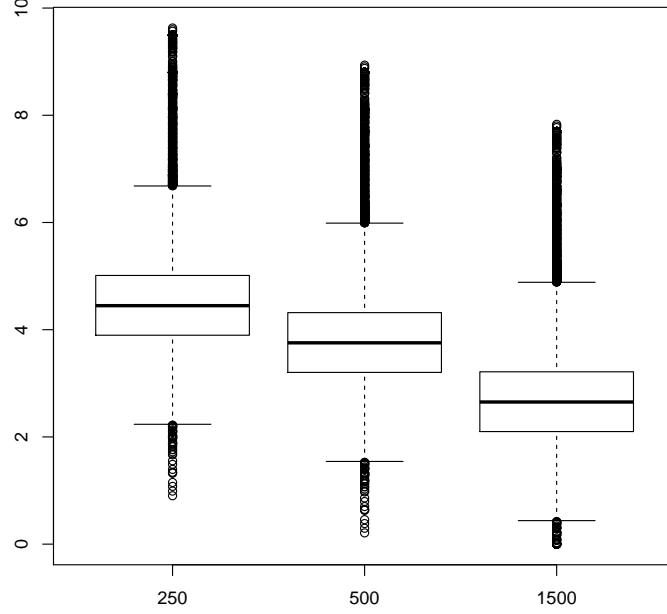


Figure 5: Boxplot of $\log(1/\pi_k)$ for different sample sizes, $n = 250, 500$ and 1500 .

A Proofs

Throughout the proofs we use the letter C to denote a generic constant whose value may vary from place to place. Let us also define $\Delta_{kl} = \pi_{kl} - \pi_k \pi_l$ and $\Delta_{kk} = \pi_k(1 - \pi_k)$.

A.1 Some useful lemmas

Lemma A.1. *Assume (A4) hold. There is a constant ζ_1 such that*

$$\frac{1}{N} \sum_{k \in U} |\phi_{k,k}(t, r)|^2 \leq \zeta_1 |t - r|^{2\beta}$$

where $\phi_{k,k}(t, r) = Y_k(t)Y_k(t) - Y_k(r)Y_k(r)$.

Proof. We have

$$\begin{aligned} \frac{1}{N} \sum_{k \in U} |\phi_{k,k}(t, r)|^2 &\leq \frac{2}{N} \left\{ \sum_U |Y_k(t) - Y_k(r)|^2 |Y_k(t)|^2 + \sum_U |Y_k(t) - Y_k(r)|^2 |Y_k(r)|^2 \right\} \\ &\leq 2 \left[\left(\frac{1}{N} \sum_U |Y_k(t) - Y_k(r)|^4 \right)^{1/2} \left(\frac{1}{N} \sum_U |Y_k(t)|^4 \right)^{1/2} \right. \\ &\quad \left. + \left(\frac{1}{N} \sum_U |Y_k(t) - Y_k(r)|^4 \right)^{1/2} \left(\frac{1}{N} \sum_U |Y_k(r)|^4 \right)^{1/2} \right]. \end{aligned}$$

Under assumption **(A4)**, we get that, for some constant ζ_1

$$\frac{1}{N} \sum_{k \in U} |\phi_{k,k}(t, r)|^2 \leq \zeta_1 |t - r|^{2\beta}.$$

□

Lemma A.2. Assume **(A3)** hold. There is a constant ζ_2 such that

$$\left(\frac{1}{N^2} \sum_{k \in U} \sum_{l \in U} |\phi_{k,l}(t, r)| \right)^2 \leq \zeta_2 |t - r|^{2\beta}$$

where $\phi_{k,l}(t, r) = Y_k(t)Y_l(t) - Y_k(r)Y_l(r)$.

Proof. The demonstration is similar to the proof of Lemma A.1 and is thus omitted. □

Lemma A.3. Assume **(A1)** and **(A2)** hold.

$$\mathbb{E}_p((\hat{d}(\pi) - d(\pi))^2) \leq \left(\frac{1}{\lambda} + \frac{\max_{k \neq l} |\Delta_{kl}|}{\lambda^2} d(\pi) \right) d(\pi).$$

Proof. Under assumptions **(A1)** and **(A2)**,

$$\begin{aligned} \mathbb{E}_p((\hat{d}(\pi) - d(\pi))^2) &= \sum_U \sum_U \frac{\pi_{kl} - \pi_k \pi_l}{\pi_k \pi_l} \pi_k (1 - \pi_k) \pi_l (1 - \pi_l) \\ &\leq \frac{1}{\lambda} \sum_U \pi_k^2 (1 - \pi_k)^2 + \frac{\max_{k \neq l} |\Delta_{kl}|}{\lambda^2} \left(\sum_U \pi_k (1 - \pi_k) \right)^2 \\ &\leq \left(\frac{1}{\lambda} + \frac{\max_{k \neq l} |\Delta_{kl}|}{\lambda^2} d(\pi) \right) d(\pi). \end{aligned}$$

□

A.2 Proof of proposition 3.1

We first consider the case of the rejective sampling $p_{rej}(s)$ and show that **A5** is true if $d(\pi_N)$ tends to infinity. By Theorem 1 in Boistard et al. (2012) and hypothesis **A2**, we have

$$\mathbb{E}_p(\mathbb{1}_{k_1 k_2 l_1 l_2}) - \pi_{k_1} \pi_{k_2} \pi_{l_1} \pi_{l_2} = O(d(\pi)^{-1})$$

uniformly for $(k_1, l_1, k_2, l_2) \in D_{4,N}$. Since $\pi_{k_1} \pi_{k_2} - \pi_{k_1 k_2} = O(d(\pi)^{-1})$ and $\pi_{l_1} \pi_{l_2} - \pi_{l_1 l_2} = O(d(\pi)^{-1})$ uniformly for $(k_1, l_1, k_2, l_2) \in D_{4,N}$, we directly obtain that, for rejective sampling

$$\max_{(k_1, l_1, k_2, l_2) \in D_{4,N}} |\mathbb{E}_p[(\mathbb{1}_{k_1 l_1} - \pi_{k_1} \pi_{l_1})(\mathbb{1}_{k_2 l_2} - \pi_{k_2} \pi_{l_2})]| \leq \frac{C}{d(\pi)},$$

for some constant C .

If we consider now a different sampling design $p_N(s)$, we have with Pinsker inequality (see Theorem 6.1 in Kemperman (1969)) and the property of the total variation distance,

$$\sup_{A \in \mathcal{A}_N} |p_N(A) - p_{rej}(A)| \leq \sqrt{K(p_N, p_{rej})/2}$$

where \mathcal{A}_N is the set of all partitions of U_N . Considering the particular cases $A = \{(k_1, l_1, k_2, l_2) \in D_{4,N}\}$, and denoting by $\pi_{k_1 k_2 l_1 l_2} = p_N(A)$ and by $\pi_{k_1 k_2 l_1 l_2}^{rej} = p_{rej}(A)$, we directly get that

$$\sup_{(k_1, l_1, k_2, l_2) \in D_{4,N}} \left| \pi_{k_1 k_2 l_1 l_2} - \pi_{k_1 k_2 l_1 l_2}^{rej} \right| \leq \sqrt{K(p_N, p_{rej})/2}$$

and the proof is complete.

A.3 Proof of Proposition 3.2 (consistency of the covariance and the variance functions)

The proof follows the same steps as in Cardot et al. (2012c). We show first that for all $t, r \in [0, T]$, the estimator of the covariance function $\hat{\gamma}_{H,d}(r, t)$ is consistent for $\gamma_p(r, t)$ and then, that the random variable $n(\hat{\gamma}_{H,d}(t, t) - \gamma_p(t, t))$ converges in distribution to zero in the space $C([0, T])$. By definition of the convergence in distribution in $C([0, T])$ and the boundedness and continuity of the sup functional, we then directly obtain the announced result. As in Cardot et al. (2012c), in order to obtain the convergence in distribution of $n(\hat{\gamma}_{H,d}(t, t) - \gamma_p(t, t))$, we first show the pointwise convergence, which clearly implies the convergence of all finite linear combinations, and then check that the sequence is tight.

Step 1. Pointwise convergence

We want to show, that for each $(t, r) \in [0, T]^2$, we have

$$n\mathbb{E}_p \{ |\hat{\gamma}_{H,d}(r, t) - \gamma_p(r, t)| \} \rightarrow 0, \quad \text{when } N \rightarrow \infty.$$

Let us decompose

$$n(\hat{\gamma}_{H,d}(r, t) - \gamma_p(r, t)) = n(\hat{\gamma}_{H,d}(r, t) - \hat{\gamma}_H(r, t)) + n(\hat{\gamma}_H(r, t) - \gamma_p(r, t))$$

and study separately the interpolation and the estimation errors.

Interpolation error

Let us suppose that $t \in [t_i, t_{i+1}[$, $r \in [t_{i'}, t_{i'+1}[$ and bound

$$\begin{aligned} n|\hat{\gamma}_{H,d}(r, t) - \hat{\gamma}_H(r, t)| &\leq \frac{n}{N^2} \frac{\hat{d}(\pi)}{d(\pi)} \sum_U \frac{1 - \pi_k}{\pi_k^2} |Y_{k,d}(t)Y_{k,d}(r) - Y_k(t)Y_k(r)| \\ &\quad + \frac{n}{N^2} \frac{1}{d(\pi)} \sum_U \sum_U \frac{1 - \pi_k}{\pi_k} \frac{1 - \pi_l}{\pi_l} |Y_{k,d}(t)Y_{l,d}(r) - Y_k(t)Y_l(r)|. \end{aligned}$$

Let us define $M = \max_{\pi_k \neq 1} \pi_k$. Noting that $\hat{d}(\pi) \leq \frac{1}{\lambda} d(\pi)$, $\frac{1}{d(\pi)} \leq \frac{1}{N\lambda(1-M)}$ and

$$\begin{aligned} |Y_{k,d}(t)Y_{l,d}(r) - Y_k(t)Y_l(r)| &\leq |Y_k(t_i) - Y_k(t)||Y_l(t_{i'})| + |Y_l(t_{i'}) - Y_l(r)||Y_k(t)| \\ &\quad + |Y_k(t_{i+1}) - Y_k(t_i)|[|Y_l(t_{i'+1})| + 2|Y_l(t_{i'})|] + |Y_k(t_i)||Y_l(t_{i'+1}) - Y_l(t_{i'})| \end{aligned}$$

we can bound, under assumptions **(A1)**-(**A4**),

$$\begin{aligned} n|\hat{\gamma}_{H,d}(r, t) - \hat{\gamma}_H(r, t)| &\leq \left(\frac{n}{N} \frac{1}{\lambda} + \frac{n}{d(\pi)} \right) \frac{C}{\lambda^2} [|t_i - t|^\beta + |t_{i'} - r|^\beta + |t_{i'+1} - t_{i'}|^\beta + 3|t_{i+1} - t_i|^\beta] \\ &\leq 2 \left(\frac{n}{N} \frac{1}{\lambda^3} + \frac{n}{d(\pi)} \right) \frac{C}{\lambda^2} [|t_{i+1} - t_i|^\beta + |t_{i'+1} - t_{i'}|^\beta]. \end{aligned} \quad (17)$$

Thus, under the assumption on the grid of discretization points,

$$n|\hat{\gamma}_{H,d}(r, t) - \hat{\gamma}_H(r, t)| = o(1).$$

Consider now the following decomposition

$$|\hat{\gamma}_H(r, t) - \gamma_p(r, t)| \leq |\hat{\gamma}_H(r, t) - \gamma_H(r, t)| + |\gamma_H(r, t) - \gamma_p(r, t)| \quad (18)$$

and study separately these two types of error.

Approximation error

We first show that, for each $(r, t) \in [0, T]^2$,

$$n|\gamma_H(r, t) - \gamma_p(r, t)| = o(1).$$

By introducing approximation (5)

$$\pi_{kl} - \pi_k \pi_l = -\pi_k \pi_l \frac{(1 - \pi_k)(1 - \pi_l)}{d(\pi)} + \frac{c_{kl}}{d(\pi)} \quad (19)$$

where $\max_{k,l} |c_{kl}| \rightarrow 0$, in the covariance function (4), we get

$$\begin{aligned} \gamma_p(r, t) &= \frac{1}{2} \frac{1}{d(\pi)N^2} \sum_{k \in U} \sum_{l \in U} [\pi_k \pi_l (1 - \pi_k)(1 - \pi_l) - c_{kl}] \left(\frac{Y_k(r)}{\pi_k} - \frac{Y_l(r)}{\pi_l} \right) \left(\frac{Y_k(t)}{\pi_k} - \frac{Y_l(t)}{\pi_l} \right) \\ &= \gamma_H(r, t) - \frac{1}{2} \frac{1}{N^2} \sum_{k \in U} \sum_{l \in U} \frac{c_{kl}}{d(\pi)} \left(\frac{Y_k(r)}{\pi_k} - \frac{Y_l(r)}{\pi_l} \right) \left(\frac{Y_k(t)}{\pi_k} - \frac{Y_l(t)}{\pi_l} \right). \end{aligned}$$

Thus, we directly get with assumptions **(A1)**-(**A3**) that

$$d(\pi) |\gamma_H(r, t) - \gamma_p(r, t)| = o(1). \quad (20)$$

Sampling error

To establish the convergence of $n(\hat{\gamma}_H(r, t) - \gamma_H(r, t))$ to zero in probability as $N \rightarrow \infty$, it is enough to show that, for all $(r, t) \in [0, T]^2$

$$n^2 \mathbb{E}_p [(\hat{\gamma}_H(r, t) - \gamma_H(r, t))^2] \rightarrow 0, \quad \text{when } N \rightarrow \infty.$$

Noting that

$$\begin{aligned} n|\hat{\gamma}_H(r, t) - \gamma_H(r, t)| &\leq \frac{n}{N^2} \left| \sum_U \left(\frac{\hat{d}(\pi)}{d(\pi)} - 1 \right) \frac{\mathbb{1}_k}{\pi_k^2} (1 - \pi_k) Y_k(t) Y_k(r) \right| \\ &\quad + \frac{n}{N^2} \left| \sum_U \left(\frac{\mathbb{1}_k}{\pi_k} - 1 \right) \frac{1 - \pi_k}{\pi_k} Y_k(t) Y_k(r) \right| \\ &\quad + \frac{n}{N^2} \frac{1}{d(\pi)} \left| \sum_U \sum_U \left(\frac{\mathbb{1}_{kl}}{\pi_k \pi_l} - 1 \right) (1 - \pi_k)(1 - \pi_l) Y_k(t) Y_l(r) \right| \\ &:= B_1(r, t) + B_2(r, t) + B_3(r, t), \end{aligned} \tag{21}$$

we get

$$n^2 \mathbb{E}_p [(\hat{\gamma}_H(r, t) - \gamma_H(r, t))^2] \leq 3\mathbb{E}_p(B_1(r, t)^2) + 3\mathbb{E}_p(B_2(r, t)^2) + 3\mathbb{E}_p(B_3(r, t)^2). \tag{22}$$

Let us show now that $\mathbb{E}_p(B_1(r, t)^2) \rightarrow 0$ when $N \rightarrow \infty$. Let us define $M = \max_{\pi_k \neq 1} \pi_k$. Under assumptions **(A1)**, **(A2)** and **(A4)** and using lemma A.3 and the inequality $\frac{1}{d(\pi)} \leq \frac{1}{N\lambda(1-M)}$, we have

$$\begin{aligned} \mathbb{E}_p(B_1(r, t)^2) &\leq \frac{n^2}{d(\pi)^2} \frac{1}{N^2} \mathbb{E}_p \left[(\hat{d}(\pi) - d(\pi))^2 \right] \left[\frac{1}{\lambda^4 N} \sum_U |Y_k(t)|^2 |Y_k(r)|^2 \right] \\ &\leq \left[\frac{1}{N(1-M)} \frac{n^2}{N^2} + n \max_{k \neq l, k, l \in U} |\Delta_{kl}| \frac{n}{N} \frac{1}{N} \right] \frac{1}{\lambda^6} \left(\frac{1}{N} \sum_U |Y_k(t)|^4 \right)^{1/2} \left(\frac{1}{N} \sum_U |Y_k(r)|^4 \right)^{1/2} \\ &\leq \frac{1}{N} C \end{aligned}$$

so that $\mathbb{E}_p(B_1(r, t)^2) \rightarrow 0$ when $N \rightarrow \infty$.

Under assumptions **(A1)**, **(A2)** and **(A4)**, we can bound

$$\begin{aligned} \mathbb{E}_p(B_2(r, t)^2) &\leq \frac{n^2}{N^4} \sum_U \sum_U \frac{|\Delta_{kl}|}{\pi_k \pi_l} \frac{1 - \pi_k}{\pi_k} \frac{1 - \pi_l}{\pi_l} |Y_k(t) Y_k(r) Y_l(t) Y_l(r)| \\ &\leq \frac{1}{\lambda^3} \frac{1}{N} \left(\frac{n^2}{N^2} + \frac{n^2 \max_{k \neq l, k, l \in U} |\Delta_{kl}|}{N\lambda} \right) \left(\frac{1}{N} \sum_U |Y_k(t)|^4 \right)^{1/2} \left(\frac{1}{N} \sum_U |Y_k(r)|^4 \right)^{1/2} \\ &\leq \frac{1}{N} C \end{aligned}$$

so that $\mathbb{E}_p(B_2(r, t)^2) \rightarrow 0$ when $N \rightarrow \infty$. For the third term, we have

$$\begin{aligned}
\mathbb{E}_p(B_3(r, t)^2) &= n^2 \mathbb{E}_p \left[\frac{1}{N^4} \frac{1}{d(\pi)^2} \sum_{k, l \in U} \sum_{k', l' \in U} \left(\frac{\mathbb{1}_{kl}}{\pi_k \pi_l} - 1 \right) \left(\frac{\mathbb{1}_{k'l'}}{\pi_{k'} \pi_{l'}} - 1 \right) \right. \\
&\quad \left. \cdot (1 - \pi_k)(1 - \pi_l)(1 - \pi_{k'})(1 - \pi_{l'}) Y_k(t) Y_l(r) Y_{k'}(t) Y_{l'}(r) \right] \\
&\leq \frac{n^2}{N^4} \frac{1}{d(\pi)^2} \sum_{k \in U} \sum_{k' \in U} \left| \mathbb{E}_p \left[\left(\frac{\mathbb{1}_k}{\pi_k^2} - 1 \right) \left(\frac{\mathbb{1}_{k'}}{\pi_{k'}^2} - 1 \right) \right] \right| |Y_k(t)| |Y_k(r)| |Y_{k'}(t)| |Y_{k'}(r)| \\
&\quad + \frac{2n^2}{N^4} \frac{1}{d(\pi)^2} \sum_{k \in U} \sum_{k' \neq l' \in U} \left| \mathbb{E}_p \left[\left(\frac{\mathbb{1}_k}{\pi_k^2} - 1 \right) \left(\frac{\mathbb{1}_{k'l'}}{\pi_{k'} \pi_{l'}} - 1 \right) \right] \right| |Y_k(t)| |Y_k(r)| |Y_{k'}(t)| |Y_{l'}(r)| \\
&\quad + \frac{n^2}{N^4} \frac{1}{d(\pi)^2} \sum_{k \neq l \in U} \sum_{k' \neq l' \in U} \left| \mathbb{E}_p \left[\left(\frac{\mathbb{1}_{kl}}{\pi_k \pi_l} - 1 \right) \left(\frac{\mathbb{1}_{k'l'}}{\pi_{k'} \pi_{l'}} - 1 \right) \right] \right| |Y_k(t)| |Y_l(r)| |Y_{k'}(t)| |Y_{l'}(r)| \\
&:= v_1 + v_2 + v_3
\end{aligned}$$

Under assumptions **(A1)**, **(A2)** and **(A4)** and the inequality $\pi_{kl} \leq \pi_k \pi_l$, we get

$$\begin{aligned}
v_1 &\leq \frac{n^2}{N^4 d(\pi)^2} \sum_{k \in U} \left| \mathbb{E}_p \left[\frac{\mathbb{1}_k}{\pi_k^4} - 2 \frac{\mathbb{1}_k}{\pi_k^2} + 1 \right] \right| |Y_k(t)|^2 |Y_k(r)|^2 \\
&\quad + \frac{n^2}{N^4 d(\pi)^2} \sum_{k \in U} \sum_{k' \neq k \in U} \left| \mathbb{E}_p \left[\frac{\mathbb{1}_{kk'}}{\pi_k^2 \pi_{k'}^2} - \frac{\mathbb{1}_k}{\pi_k^2} - \frac{\mathbb{1}_{k'}}{\pi_{k'}^2} + 1 \right] \right| |Y_k(t)| |Y_k(r)| |Y_{k'}(t)| |Y_{k'}(r)| \\
&\leq \frac{n^2}{N^2 d(\pi)^2} \left[\frac{1}{N} \left(\frac{1}{\lambda^3} + \frac{2}{\lambda} + 1 \right) + \left(\frac{1}{\lambda^2} + \frac{2}{\lambda} + 1 \right) \right] \left(\frac{1}{N} \sum_{k \in U} Y_k(t)^4 \right)^{1/2} \left(\frac{1}{N} \sum_{k \in U} Y_k(r)^4 \right)^{1/2}.
\end{aligned} \tag{23}$$

Since $d(\pi) \rightarrow \infty$ when $N \rightarrow \infty$ we have that $v_1 \rightarrow 0$. Under assumptions **(A1)**, **(A2)**, **(A4)** and **(A5)**

$$\begin{aligned}
v_3 &\leq \frac{C}{N} + \frac{n^2}{N^4 d(\pi)^2} \frac{1}{\lambda^4} \max_{(k, l, k', l') \in D_{4, N}} |\mathbb{E}_p [(\mathbb{1}_{kl} - \pi_k \pi_l)(\mathbb{1}_{k'l'} - \pi_{k'} \pi_{l'})]| \sum_{(k, l, k', l') \in D_{4, N}} |Y_k(t)| |Y_l(r)| |Y_{k'}(t)| |Y_{l'}(r)| \\
&\leq \frac{C}{N} + \frac{n^2}{\lambda^4 d^2(\pi)} \max_{(k, l, k', l') \in D_{4, N}} |\mathbb{E}_p [(\mathbb{1}_{kl} - \pi_k \pi_l)(\mathbb{1}_{k'l'} - \pi_{k'} \pi_{l'})]| \left(\frac{1}{N} \sum_{k \in U} |Y_k(t)|^2 \right) \left(\frac{1}{N} \sum_{l \in U} |Y_l(r)|^2 \right).
\end{aligned} \tag{24}$$

Hence $v_3 \rightarrow 0$ when $N \rightarrow \infty$. By the Cauchy-Schwarz inequality, we have $v_2 \rightarrow 0$ when $N \rightarrow \infty$. Finally, we have for all $(r, t) \in [0, T]^2$,

$$n \mathbb{E}_p(|\hat{\gamma}_H(r, t) - \gamma_p(r, t)|) \rightarrow 0, \quad \text{when } N \rightarrow \infty. \tag{25}$$

and consequently,

$$n \mathbb{E}_p \{ |\hat{\gamma}_{H, d}(r, t) - \gamma_p(r, t)| \} \rightarrow 0, \quad \text{when } N \rightarrow \infty.$$

Step 2. Tightness

To check the tightness of $n(\hat{\gamma}_H(t, t) - \gamma_H(t, t))$ in $C[0, T]$, we use the Theorem 12.3 from Billingsley (1968). Since the pointwise consistency of $n(\hat{\gamma}_H - \gamma_H)$ implies that $n(\hat{\gamma}_H(0, 0) - \gamma_H(0, 0))$ is tight, to get the announced result, it remains to study the increments of $n(\hat{\gamma}_H - \gamma_H)$ between two instants t and r . Considering

$$d_\gamma^2(t, r) = n^2 \mathbb{E}_p(|\hat{\gamma}_H(t, t) - \gamma_H(t, t) - \hat{\gamma}_H(r, r) + \gamma_H(r, r)|^2)$$

we only need to prove that

$$d_\gamma^2(t, r) \leq C|t - r|^{2\beta},$$

for some positive constant C and all $(r, t) \in [0, T]^2$. Since $\beta > 1/2$, the above inequality implies that the sequence $n(\hat{\gamma}_H - \gamma_H)$ is tight in $C([0, T])$.

Using (21), we can decompose $d_\gamma^2(t, r)$ into 3 parts,

$$\begin{aligned} d_\gamma^2(t, r) &\leq 3 \left(\mathbb{E}_p([B_1(t, t) - B_1(r, r)]^2) + \mathbb{E}_p([B_2(t, t) - B_2(r, r)]^2) \right. \\ &\quad \left. + \mathbb{E}_p([B_3(t, t) - B_3(r, r)]^2) \right) \\ &:= 3(d_{B_1}^2 + d_{B_2}^2 + d_{B_3}^2). \end{aligned}$$

Using Lemma A.1 and Assumptions **(A1)**-(**A2**), we obtain

$$\begin{aligned} d_{B_1}^2 &= \frac{n^2}{N^4} \mathbb{E}_p \left[\left(\sum_k \frac{\hat{d}(\pi) - d(\pi)}{d(\pi)} \frac{\mathbb{1}_k}{\pi_k} \frac{1 - \pi_k}{\pi_k} \phi_{k,k}(t, r) \right)^2 \right] \\ &\leq \frac{n^2}{N^2} \left(\frac{1}{\lambda} - 1 \right)^2 \frac{1}{\lambda^4} \left(\frac{1}{N} \sum_U \phi_{k,k}(t, r) \right)^2 \\ &\leq \frac{n^2}{N^2} \left(\frac{1}{\lambda} - 1 \right)^2 \frac{1}{\lambda^4} \zeta_1 |t - r|^{2\beta} \\ &\leq C|t - r|^{2\beta}. \end{aligned} \tag{26}$$

With assumptions **(A1)**, **(A2)** and **(A6)** and Lemma A.1, we get

$$\begin{aligned} d_{B_2}^4 &\leq \frac{1}{\lambda^3} \left(\frac{n^2}{N^3} + \frac{n^2 \max_{k \neq l} |\Delta_{kl}|}{N^2 \lambda} \right) \frac{1}{N} \sum_U |\phi_{k,k}(t, r)|^2 \\ &\leq C|t - r|^{2\beta}. \end{aligned} \tag{27}$$

Finally,

$$\begin{aligned}
d_{B_3}^4 &\leq \frac{n^2}{N^4} \frac{1}{d(\pi)^2} \sum_k \sum_{k'} \left| \mathbb{E}_p \left[\left(\frac{\mathbb{1}_k}{\pi_k^2} - 1 \right) \left(\frac{\mathbb{1}_{k'}}{\pi_{k'}^2} - 1 \right) \right] \right| |\phi_{k,k}(t, r)| |\phi_{k',k'}(t, r)| \\
&\quad + \frac{2n^2}{N^4} \frac{1}{d(\pi)^2} \sum_k \sum_{k' \neq l'} \left| \mathbb{E}_p \left[\left(\frac{\mathbb{1}_k}{\pi_k^2} - 1 \right) \left(\frac{\mathbb{1}_{k'l'}}{\pi_{k'} \pi_{l'}} - 1 \right) \right] \right| |\phi_{k,k}(t, r)| |\phi_{k',l'}(t, r)| \\
&\quad + \frac{n^2}{N^4} \frac{1}{d(\pi)^2} \sum_{k \neq l} \sum_{k' \neq l'} \left| \mathbb{E}_p \left[\left(\frac{\mathbb{1}_{kl}}{\pi_k \pi_l} - 1 \right) \left(\frac{\mathbb{1}_{k'l'}}{\pi_{k'} \pi_{l'}} - 1 \right) \right] \right| |\phi_{k,l}(t, r)| |\phi_{k',l'}(t, r)| \\
&:= b_1 + b_2 + b_3.
\end{aligned}$$

Thanks to Lemma A.1 and under assumptions **(A1)**, **(A2)** and **(A4)** we get

$$\begin{aligned}
b_1 &\leq \frac{n^2}{N^2} \frac{1}{d(\pi)^2} \left[\frac{1}{N} \left(\frac{1}{\lambda^3} + \frac{2}{\lambda} + 1 \right) + \left(\frac{1}{\lambda^2} + \frac{2}{\lambda} + 1 \right) \right] \frac{1}{N} \sum_k |\phi_{k,k}(t, r)|^2 \\
&\leq C|t - r|^{2\beta}.
\end{aligned} \tag{28}$$

Under assumptions **(A1)**, **(A2)**, **(A4)** and **(A5)** and using Lemma A.2, we have

$$\begin{aligned}
b_3 &\leq \frac{C|t - r|^{2\beta}}{N} + \frac{n^2}{d^2(\pi)\lambda^4} \max_{(k,l,k',l') \in D_{4,N}} |\mathbb{E}_p [(\mathbb{1}_{kl} - \pi_k \pi_l)(\mathbb{1}_{k'l'} - \pi_{k'} \pi_{l'})]| \left(\frac{1}{N^2} \sum_{k,l} |\phi_{k,l}(t, r)| \right)^2 \\
&\leq C|t - r|^{2\beta}.
\end{aligned} \tag{29}$$

Using the Cauchy-Schwarz inequality together with bounds (28) and (29), we get that $b_2 \leq C|t - r|^{2\beta}$ so that

$$d_{B_3}^2 \leq C|t - r|^{2\beta}. \tag{30}$$

Finally, we deduce, with inequalities (26), (27) and (30) that

$$d_\gamma^2(r, t) \leq C|t - r|^{2\beta}. \tag{31}$$

A.4 Proof of Proposition 3.3

Under assumptions **(A1)** and **(A2)**, it is clear that $\hat{d}(\pi)/d(\pi) = 1 + o_p(1)$. The pointwise convergence of $n\hat{\gamma}_{H,d}^*(r, t)$ is then a direct consequence of Proposition 3.2 and the fact that $\hat{\gamma}_{H,d}^*(r, t) = \frac{d(\pi)}{\hat{d}(\pi)} \hat{\gamma}_{H,d}(r, t)$. Furthermore, we may write

$$n(\hat{\gamma}_{H,d}^* - \gamma_H) = n \frac{d(\pi)}{\hat{d}(\pi)} (\hat{\gamma}_{H,d} - \gamma_H) + n \left(\frac{d(\pi)}{\hat{d}(\pi)} - 1 \right) \gamma_H.$$

By Slutsky's theorem, the first term at the righthand-side of previous equation converges in distribution to zero in $C([0, T])$ while the second term goes to zero in probability since $\sup_{(r,t) \in [0,T]^2} |n\gamma_H(r, t)| < \infty$ and $\frac{d(\pi)}{\hat{d}(\pi)} - 1 = o_p(1)$. Hence, the sequence $n(\hat{\gamma}_{H,d}^* - \gamma_H)$ converges in distribution to zero in $C([0, T])$.

A.5 Proof of Proposition 3.4

We first note that the interpolation error, bounded in (17), satisfies

$$n^{3/2} |\hat{\gamma}_{H,d}(r, t) - \hat{\gamma}_H(r, t)| = O(1) \quad (32)$$

provided that $\lim_{N \rightarrow \infty} \max_{i=\{1, \dots, d_N-1\}} |t_{i+1} - t_i|^{2\beta} = O(n^{-1})$. We then use the fact (see Theorem 1 in Boistard et al. (2012)) that for rejective sampling the terms c_{kl} defined in (19) satisfy, for some constant C ,

$$\max_{k,l} |c_{kl}| \leq Cd(\pi)^{-1}.$$

Thus, bound (20) is now

$$d(\pi)^2 |\gamma_H(r, t) - \gamma_p(r, t)| = O(1). \quad (33)$$

If we examine now the sampling error, we can check that the terms B_1 and B_2 are of order n^{-1} . Concerning the term B_3 , it is bounded by the sum $v_1 + v_2 + v_3$ with $v_1 = O(d^{-2}(\pi))$ and $v_2 \leq \sqrt{v_1 v_3}$. Thanks to Proposition 3.1, we get that the term v_3 satisfies $v_3 = O(d^{-1}(\pi))$ and consequently, $\mathbb{E}_p(B_3(r, t)^2) = O(n^{-1})$. Thus,

$$n^2 \mathbb{E}_p [(\hat{\gamma}_H(r, t) - \gamma_H(r, t))^2] = O(n^{-1})$$

and the proof is complete.

References

- Beaumont, J.-F. and Rivest, L.-P. (2009). Dealing with outliers in survey data. In *Sample surveys: design, methods and applications*, volume 29 of *Handbook of Statist.*, pages 247–279. Elsevier/North-Holland, Amsterdam.
- Berger, Y. G. (1998a). Rate of convergence for asymptotic variance of the horvitz-thompson estimator. *Journal of Statistical Planning and Inference*, 74:149–168.
- Berger, Y. G. (1998b). Rate of convergence to normal distribution for the horvitz-thompson estimator. *Journal of Statistical Planning and Inference*, 67:209–226.
- Billingsley, P. (1968). *Convergence of Probability Measures*. John Wiley and Sons.
- Boistard, H., Lopuhaä, H. P., and Ruiz-Gazen, A. (2012). Approximation of rejective sampling inclusion probabilities and application to higher order correlation. Technical report, Arxiv 1207.5654.

- Bondesson, L. (2010). Conditional and restricted Pareto sampling: two new methods for unequal probability sampling. *Scand. J. Stat.*, 37(3):514–530.
- Bondesson, L., Traat, I., and Lundqvist, A. (2006). Pareto sampling versus Sampford and conditional Poisson sampling. *Scand. J. Statist.*, 33(4):699–720.
- Brewer, K. R. W. and Hanif, M. (1983). *Sampling with unequal probabilities*, volume 15 of *Lecture Notes in Statistics*. Springer-Verlag, New York.
- Cardot, H., Chaouch, M., Goga, C., and Labruère, C. (2010). Properties of design-based functional principal components analysis. *J. of Statistical Planning and Inference*, 140:75–91.
- Cardot, H., Degras, D., and Josserand, E. (2012a). Confidence bands for Horvitz-Thompson estimators using sampled noisy functional data. *Bernoulli*, to appear.
- Cardot, H., Dessertaine, A., Goga, C., Josserand, E., and Lardin, P. (2012b). Comparaison de différents plans de sondage et construction de bandes de confiance pour l’estimation de la moyenne de données fonctionnelles : une illustration sur la consommation électrique. *Survey Methodology / Technique d’enquêtes*, to appear.
- Cardot, H., Goga, C., and Lardin, P. (2012c). Uniform convergence and asymptotic confidence bands for model-assisted estimators of the mean of sampled functional data. Technical report, arXiv:1204.6382.
- Cardot, H. and Josserand, E. (2011). Horvitz-Thompson estimators for functional data: asymptotic confidence bands and optimal allocation for stratified sampling. *Biometrika*, 98:107–118.
- Chauvet, G. (2007). *Méthodes de bootstrap en population finie*. PhD thesis, Université de Rennes II.
- Chauvet, G. and Tillé, Y. (2006). A fast algorithm of balanced sampling. *Computational Statistics*, 21:53–61.
- Chen, X.-H., Dempster, A. P., and Liu, J. S. (1994). Weighting finite population sampling to maximise entropy. *Biometrika*, 81:457–469.
- Degras, D. A. (2011). Simultaneous confidence bands for nonparametric regression with functional data. *Statist. Sinica*, 21:1735–1765.
- Deville, J.-C. and Tillé, Y. (2004). Efficient balanced sampling: the cube algorithm. *Biometrika*, 91:893–912.

- Deville, J.-C. and Tillé, Y. (2005). Variance approximation under balanced sampling. *J. Statist. Plann. Inference*, 128:569–591.
- Erdős, P. and Rényi, A. (1959). On the central limit theorem for samples from a finite population. *Publ. Math. Inst. Hungar. Acad.Sci.*, 4:49–61.
- Fuller, W. A. (2009). Some design properties of a rejective sampling procedure. *Biometrika*, 96(4):933–944.
- Hájek, J. (1960). Limiting distributions in simple random sampling from a finite population. *Publ. Math. Inst. Hungar. Acad.Sci.*, 77:361–374.
- Hájek, J. (1964). Asymptotic theory of rejective sampling with varying probabilities from a finite population. *Annals of Mathematical Statistics*, 35:1491–1523.
- Hájek, J. (1981). *Sampling from a finite population*. Statistics: Textbooks and Monographs. Marcel Dekker, New York.
- Kemperman, J. H. B. (1969). On the optimum rate of transmitting information. *Ann. Math. Statist.*, 40:2156–2177.
- Lundqvist, A. (2007). On the distance between some π ps sampling designs. *Acta Appl. Math.*, 97(1-3):79–97.
- Matei, A. and Tillé, Y. (2005). Evaluation of variance approximations and estimators in maximum entropy sampling with unequal probability and fixed sample size. *Journal of Official Statistics*, 21(4):543–570.
- Ramsay, J. O. and Silverman, B. W. (2005). *Functional data analysis*. Springer-Verlag, New York, second edition.
- Sen, A. (1953). On the estimate of the variance in sampling with varying probabilities. *Journal of the Indian Society of Agricultural Statistics*, 5:119–127.
- Tillé, Y. (2006). *Sampling algorithms*. Springer Series in Statistics. Springer, New York.
- Víšek, J. Á. (1979). Asymptotic distribution of simple estimate for rejective, Sampford and successive sampling. In *Contributions to statistics*, pages 263–275. Reidel, Dordrecht.
- Yates, F. and Grundy, P. (1953). Selection without replacement from within strata with probability proportional to size. *J. Royal Statist. Soc., B*, 15:235–261.